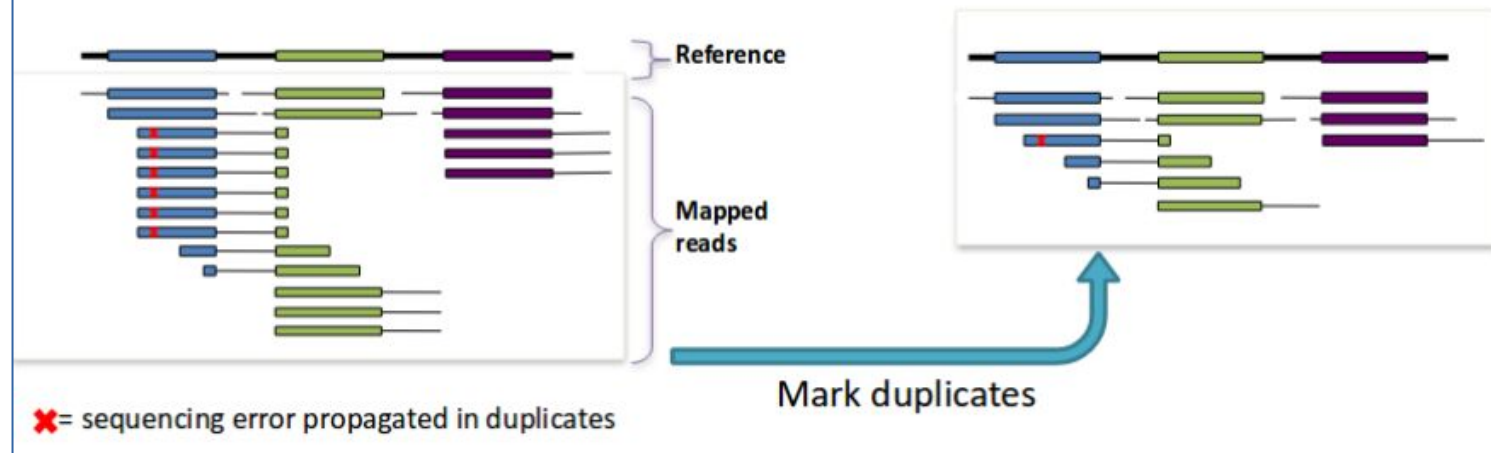# Module 1/5: Analyses ADN

- NGS Introduction

- Reads Quality Control

- Reads Cleaning

- Aligning reads on reference          → *Hélène Touzet*

- Alignment parameters          → *Hélène Touzet*

- Reads duplicates

→ **Practical #3**

# Cleaning duplicated reads



## Why mark duplicates?

- Duplicates = sets of reads pairs with same unclipped alignment start and unclipped alignment end
- Suspected to be **non-independent measurements** of a sequence
  - Sampled from the exact same template of DNA
  - Violates assumptions of variant calling
- Errors in sample/library prep will get propagated to *all* the duplicates
  - Just pick the "best" copy – mitigates the effects of errors

Reference

Mapped reads

Mark duplicates

✖ = sequencing error propagated in duplicates

*Source : GATK Marking duplicates*
*https://software.broadinstitute.org/gatk/events/slides/1511/Presentations/GATKwh9-3-Marking_duplicates.pdf*

# Picard / MarkDuplicate

Additional information about Picard tools is available from Picard web site at *http://broadinstitute.github.io/picard/*



|  | BWA | Bowtie2 |
|---|---|---|
| UNPAIRED_READS_EXAMINED | 18 | 19 |
| READ_PAIRS_EXAMINED | 4043 | 4044 |
| SECONDARY_OR_SUPPLEMENTARY_RDS | 0 | 3 |
| UNMAPPED_READS | 22 | 19 |
| UNPAIRED_READ_DUPLICATES | 0 | 0 |
| READ_PAIR_DUPLICATES | 12 | 12 |
| READ_PAIR_OPTICAL_DUPLICATES | 0 | 0 |
| PERCENT_DUPLICATION | 0,002962 | 0,00296 |
| ESTIMATED_LIBRARY_SIZE | 679728 | 680065 |

# Alignment count : *samtools flagstat*



## BWA

```
8129 + 0 in total (QC-passed reads + QC-failed reads)
0 + 0 secondary
3 + 0 supplementary
24 + 0 duplicates
8110 + 0 mapped (99.77%:-nan%)
8126 + 0 paired in sequencing
4063 + 0 read1
4063 + 0 read2
7980 + 0 properly paired (98.20%:-nan%)
8088 + 0 with itself and mate mapped
19 + 0 singletons (0.23%:-nan%)
0 + 0 with mate mapped to a different chr
0 + 0 with mate mapped to a different chr (mapQ>=5)
```

## Bowtie 2

```
8126 + 0 in total (QC-passed reads + QC-failed reads)
0 + 0 secondary
0 + 0 supplementary
24 + 0 duplicates
8104 + 0 mapped (99.73%:-nan%)
8126 + 0 paired in sequencing
4063 + 0 read1
4063 + 0 read2
8074 + 0 properly paired (99.36%:-nan%)
8086 + 0 with itself and mate mapped
18 + 0 singletons (0.22%:-nan%)
0 + 0 with mate mapped to a different chr
0 + 0 with mate mapped to a different chr (mapQ>=5)
```

# Coverage and deepth of coverage



Source : Élodie Girard , 5ème Ecole de bioinformatique AVIESAN-IFB 2016
http://www.france-bioinformatique.fr/sites/default/files/V01_ITMO_2016_EG_from_fastq_to_mapping_1.pdf

# Computing coverage and deepth of coverage
## *DeepTools2 / plotCoverage*

# DeepTools / Plot Coverage

# Galaxy *Workflow*

- Extract *workflow* from an history
- Modify *workflow*
- Execute *workflow* on new data
- Compare results from 2 *workflows* (in 2 histories)

# Extract *Workflow* from the history of steps applied to the first sample

# Visualize workflow

# Modify workflow visualisation

# Modify some steps configuration

This WF uses 3 input files. Change box name to describe which data is required for each input : eg *Reference, Forward fastq, Reverse fastq*



You can also change any parameter for example for *trimmomatic* step.

# Enable a parameter to be set at run time

Parameters for each tool will have the predefined values set in the workflow
You can modify this to enable any parameter to be set at run time.
Modify Trimmomatic so that Adapter are set at run time



Do'nt forget to save your workflow !

# Import new data for sample HG0103

Importe files HG0103_1.fastq and HG_0103_2.fastq

# Analyze these new data
# with the same workflow

Run the workflow with these new data



**Galaxy**    Analyze Data    **Workflow**    Shared Data ▾    Visualization ▾

## Your workflows

⊕ Create new workflow    ⬆ Upload or import workflow

| Name | # of Steps |
|---|---|
| TP1_WF1_OK ▾ | 17 |
|  | 17 |

Edit
Run
Share or Download
Copy
Rename
View
Delete

**Wo** with you by others

No w ith you.

Ot

---

**Workflow: TP1_WF1_OK**    ✔ Run workflow

**History Options**

Send results to a new history

[ Yes ] No

History name

HG0103

**Step 1: Input dataset**

Forward fastq

27 : HG00103_1.fastq ▾

**Step 2: Input dataset**

Reverse fastq

28 : HG00103_2.fastq ▾

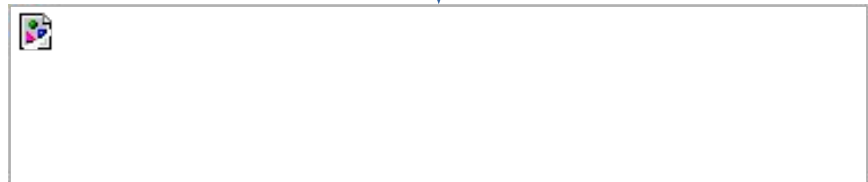**Step 3: Input dataset**

Reference

1 : GRCh37_region1.fasta ▾

**Step 4: FASTQ Groomer** convert between various FASTQ quality formats (Galaxy Version 1.0.4)

**Step 5: FASTQ Groomer** convert between various FASTQ quality formats (Galaxy Version 1.0.4)

**Step 6: FastQC** Read Quality reports (Galaxy Version 0.67)

Short read data from your current history

# Browse results

# Gather BWA alignment results
# for the 2 samples 1/2

Create a new history named results
From history TP1 : Copy HG0101_BWA_MD.bam *dataset*

# Gather BWA alignment results
# for the 2 samples 2/2

From history HG0103 : Copy *dataset*
« MarkDuplicates on data 13: MarkDuplicates BAM output »

# Visualize deepth of coverage
# for both samples

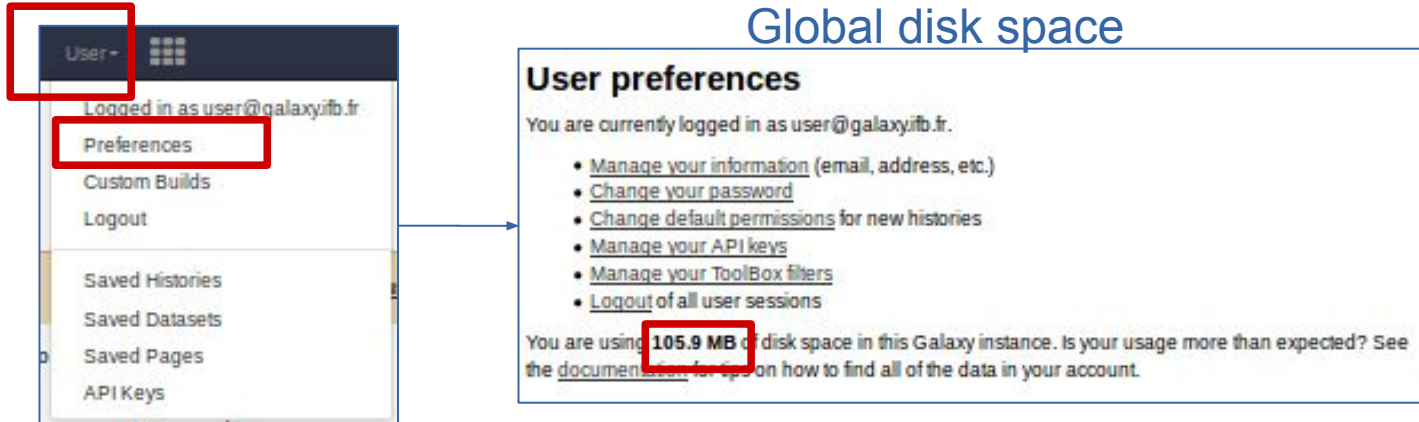Rename *datasets*
Run *plotCoverage*

# Galaxy – *Best Practices*

- Manage disk space

- Export analysis results (*datasets and histories*)

- Export / Import analysis protocoles (*workflow)*

# Manage disk space



## Global disk space
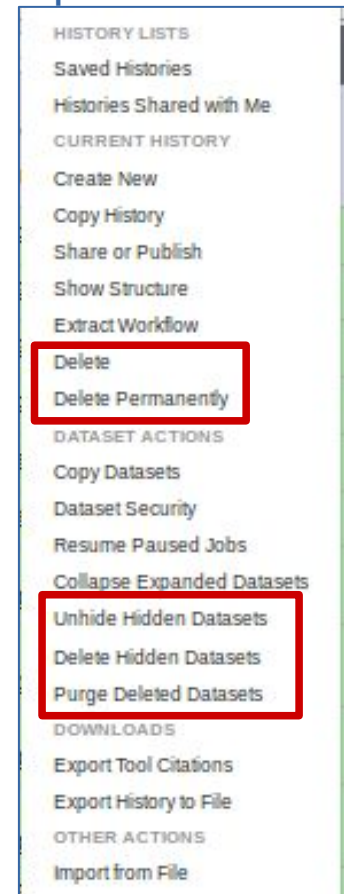
**User preferences**

You are currently logged in as user@galaxy.ifb.fr.

- Manage your information (email, address, etc.)
- Change your password
- Change default permissions for new histories
- Manage your API keys
- Manage your ToolBox filters
- Logout of all user sessions

You are using **105.9 MB** of disk space in this Galaxy instance. Is your usage more than expected? See the documentation for tips on how to find all of the data in your account.

**User-**

- Logged in as user@galaxy.ifb.fr
- Preferences
- Custom Builds
- Logout

- Saved Histories
- Saved Datasets
- Saved Pages
- API Keys

## Disk space per *dataset*

HISTORY LISTS
Saved Histories
Histories Shared with Me
CURRENT HISTORY
Create New
Copy History
Share or Publish
Show Structure
Extract Workflow
Delete
Delete Permanently
DATASET ACTIONS
Copy Datasets
Dataset Security
Resume Paused Jobs
Collapse Expanded Datasets
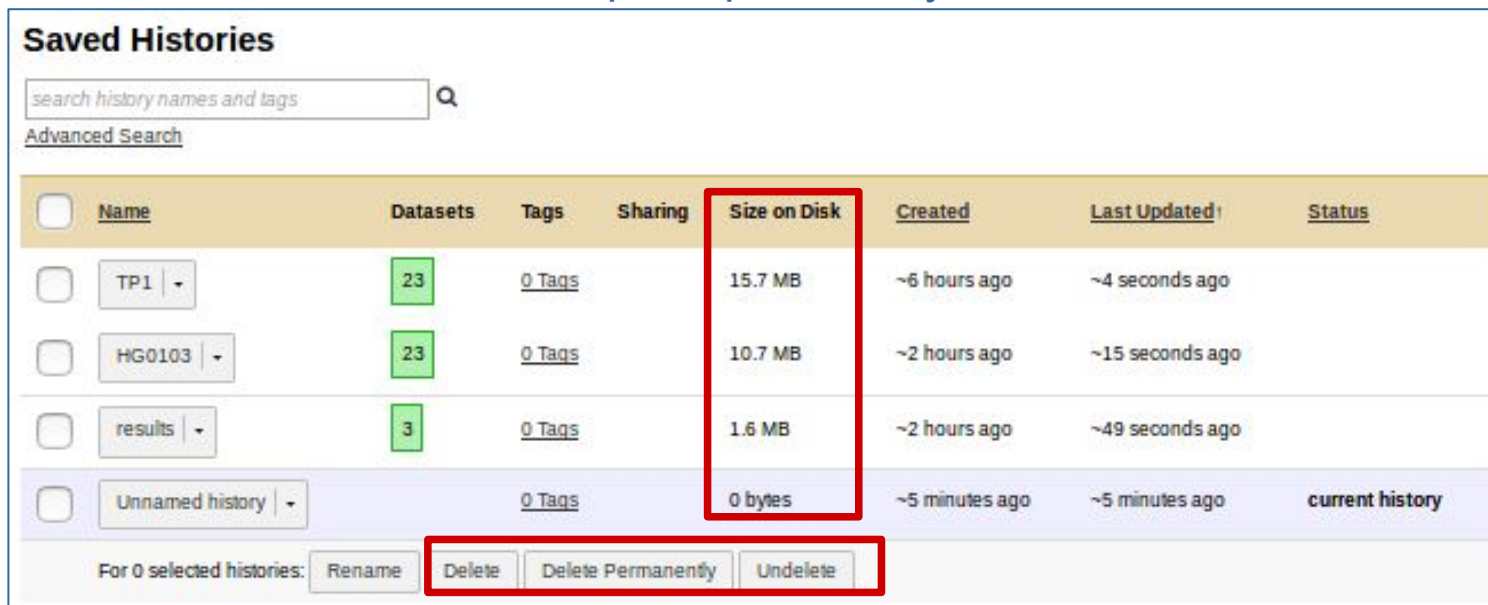Unhide Hidden Datasets
Delete Hidden Datasets
Purge Deleted Datasets
DOWNLOADS
Export Tool Citations
Export History to File
OTHER ACTIONS
Import from File

## Disk space per history

**Saved Histories**

search history names and tags 🔍

Advanced Search

| | Name | Datasets | Tags | Sharing | Size on Disk | Created | Last Updated↑ | Status |
|---|---|---|---|---|---|---|---|---|
| ☐ | TP1 ▾ | 23 | 0 Tags | | 15.7 MB | ~6 hours ago | ~4 seconds ago | |
| ☐ | HG0103 ▾ | 23 | 0 Tags | | 10.7 MB | ~2 hours ago | ~15 seconds ago | |
| ☐ | results ▾ | 3 | 0 Tags | | 1.6 MB | ~2 hours ago | ~49 seconds ago | |
| ☐ | Unnamed history ▾ | | 0 Tags | | 0 bytes | ~5 minutes ago | ~5 minutes ago | **current history** |

For 0 selected histories: Rename Delete Delete Permanently Undelete

# Export analysis results : *datasets*

# Export analysis results : *histories*

# Export / import analysis protocoles  : *workflow*

Export

**Your workflows**

**Name**

TP1_WF1_OK ▾

Edit

Run

Share or Download

Copy

Rename

View

Delete

Import